

# Basic Optimization Algorithms

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Importance of Optimization in machine learning
- How learning differs from optimization
- Challenges in neural network optimization
- **Basic Optimization Algorithms**
  1. Stochastic Gradient Descent
  2. Momentum
  3. Nesterov Momentum
- Parameter initialization strategies
- Algorithms with adaptive learning rates
- Approximate second-order methods
- Optimization strategies and meta-algorithms

# Stochastic Gradient Descent

- We have seen:
  - Gradient descent that follows the gradient of an entire training set downhill
  - This can be accelerated considerably by SGD which follows the gradient of randomly selected minibatches downhill
- SGD and its variants are the most used optimization algorithms for ML in general and deep learning in particular
  - Average gradient on a minibatch is an estimate of the gradient

# Following gradient estimate downhill

- Algorithm: SGD update at training iteration  $k$

**Require:** Learning rate  $\epsilon_k$ .

**Require:** Initial parameter  $\theta$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient estimate:  $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

    Apply update:  $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

**end while**

---

- A crucial parameter for Algorithm SGD is the learning rate  $\epsilon$ 
  - It is necessary to let the learning rate decrease
    - At iteration  $k$  it is  $\epsilon_k$

# Need for decreasing learning rate

- Batch gradient descent can use a fixed learning rate
  - Since true gradient becomes small and then 0
- SGD has a source of error
  - Random sampling of  $m$  training samples
  - Sufficient condition for SGD convergence  $\sum_{k=1}^{\infty} \epsilon_k = \infty \quad \sum_{k=1}^{\infty} \epsilon_k^2 < \infty$
  - Common to decay learning rate linearly until iteration  $\tau$ :  $\epsilon_k = (1-\alpha)\epsilon_0 + \alpha\epsilon_\tau$  with  $\alpha = k/\tau$ .
  - After iteration  $\tau$ , it is common to leave  $\epsilon$  constant

## 2. The Momentum method

- SGD is a popular optimization strategy
- But it can be slow
- Momentum method accelerates learning, when:
  - Facing high curvature
  - Small but consistent gradients
  - Noisy gradients
- Algorithm accumulates moving average of past gradients and move in that direction, while exponentially decaying

# Momentum definition

- Introduce variable  $\mathbf{v}$ , or velocity
- It is the direction and speed at which parameters move through parameter space
- Momentum in physics is mass times velocity
- The momentum algorithm assumes unit mass
- A hyperparameter  $\alpha \in [0,1)$  determines exponential decay

# Momentum update rule

- The update rule is given by

$$\begin{aligned} \mathbf{v} &\leftarrow \alpha \mathbf{v} - \varepsilon \nabla_{\theta} \left( \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}) \right) \\ \theta &\leftarrow \theta + \mathbf{v} \end{aligned}$$

- The velocity  $\mathbf{v}$  accumulates the gradient elements  $\nabla_{\theta} \left( \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}) \right)$
- The larger  $\alpha$  is relative to  $\varepsilon$ , the more previous gradients affect the current direction
- The SGD algorithm with momentum is next



# SGD algorithm with momentum

- Stochastic gradient descent with momentum

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

**Require:** Initial parameter  $\theta$ , initial velocity  $v$ .

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute gradient estimate:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

    Compute velocity update:  $v \leftarrow \alpha v - \epsilon g$

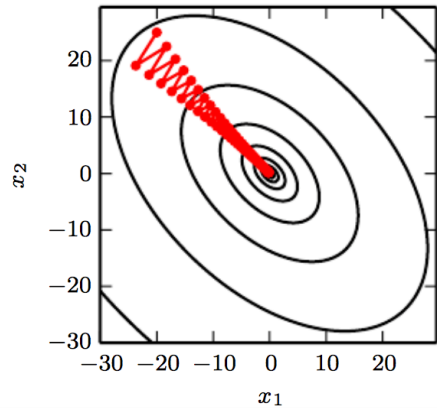
    Apply update:  $\theta \leftarrow \theta + v$

**end while**

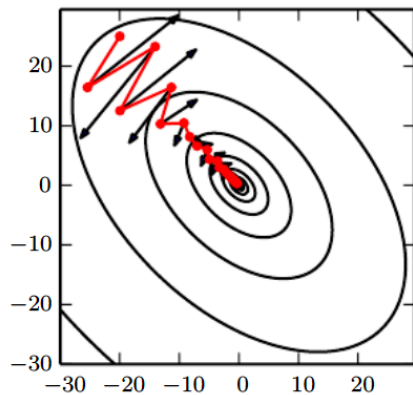
---

# Momentum

- SGD without momentum



- SGD with momentum



# Nesterov Momentum

- A variant of the momentum algorithm
- An accelerated gradient method
- It applies a correction factor to the standard method

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

**Require:** Initial parameter  $\theta$ , initial velocity  $v$ .

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding labels  $y^{(i)}$ .

    Apply interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$

    Compute gradient (at interim point):  $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$

    Compute velocity update:  $v \leftarrow \alpha v - \epsilon g$

    Apply update:  $\theta \leftarrow \theta + v$

**end while**